



# Deep Bayesian Self-Training

Fabio De Sousa Ribeiro<sup>1</sup> · Francesco Calivá<sup>1</sup> · Mark Swainson<sup>2</sup> · Kjartan Gudmundsson<sup>2</sup> · Georgios Leontidis<sup>1</sup> · Stefanos Kollias<sup>1</sup>

Received: 3 November 2018 / Accepted: 28 June 2019  
© The Author(s) 2019

## Abstract

Supervised deep learning has been highly successful in recent years, achieving state-of-the-art results in most tasks. However, with the ongoing uptake of such methods in industrial applications, the requirement for large amounts of annotated data is often a challenge. In most real-world problems, manual annotation is practically intractable due to time/labour constraints; thus, the development of automated and adaptive data annotation systems is highly sought after. In this paper, we propose both a (1) deep Bayesian self-training methodology for automatic data annotation, by leveraging predictive uncertainty estimates using variational inference and modern neural network (NN) architectures, as well as (2) a practical adaptation procedure for handling high label variability between different dataset distributions through clustering of NN latent variable representations. An experimental study on both public and private datasets is presented illustrating the superior performance of the proposed approach over standard self-training baselines, highlighting the importance of predictive uncertainty estimates in safety-critical domains.

**Keywords** Bayesian CNN · Variational inference · Self-training · Uncertainty weighting · Deep learning · Clustering · Representation learning · Adaptation

## 1 Introduction

With the advent of Big Data in industrial applications, the ability to automatically label datasets using limited supervision is increasingly sought after. In most real-world

problems, manual annotation is practically intractable due to time and labour constraints. Furthermore, recent advances in supervised deep learning have shown that training over parameterised models on large datasets significantly increases performance [1]. With that in mind—and despite the high demand for annotated data—deep learning practitioners have not yet explored or leveraged many of deep learning tools for automatic annotation systems. This is evidenced by the scarcity of existing research in the field, compared to others [2]. Automated annotation techniques typically involve semi-supervised algorithmic variants, wherein learning systems are often trained on a small initial sample of labelled data, and leverage information from unlabelled data to generalise better [3]. Well-established semi-supervised methods such as self-training [4], transfer learning [5], co-training [6], active learning [7] and tri-training [8] among others have shown to be useful for labelling in the past, but some challenges remain with regard to their scalability to high-dimensional data and their suitability to modern deep learning settings [2, 9]. Prominent recent works have explored some of these ideas in the context of modern deep models, proposing new paradigms such as co-teaching [10], active learning on

---

✉ Georgios Leontidis  
gleontidis@lincoln.ac.uk

Fabio De Sousa Ribeiro  
fdesousaribeiro@lincoln.ac.uk

Francesco Calivá  
fcaliva@lincoln.ac.uk

Mark Swainson  
mswainson@lincoln.ac.uk

Kjartan Gudmundsson  
kgudmundsson@lincoln.ac.uk

Stefanos Kollias  
skollias@lincoln.ac.uk

<sup>1</sup> MLearn Group, School of Computer Science, University of Lincoln, Lincoln, UK

<sup>2</sup> National Centre for Food Manufacturing, Holbeach Technology Park, Holbeach, UK

image data [2] and analysing deep transfer learning [11, 12] with good levels of success. Taking inspiration from these works, in this paper we primarily focus on exploring the self-training algorithm in combination with modern Bayesian deep learning methods and leverage predictive uncertainty estimates for self-labelling of high-dimensional data.

## 1.1 Background on application domain

In addition to public domain datasets, we evaluate our methods on a real-world task involving optical character verification (OCV) of real food packaging images, expanding on earlier work in [13] by reducing manual data annotation.

Incorrectly labelled food products (e.g. bearing an incorrect/illegible *use-by* date) result in product recalls and food waste, as label faults can lead to food safety incidents. Label faults are primarily attributed to human error during error-prone manual checking. Automatic approaches typically involve OCV, whereby a supervisory system holds the correct date code string and transfers it to both the printer and the vision system. The latter will then verify its read and take appropriate action. Such a system could also be used alongside other systems such as blockchain, within the food chain for food traceability [14]. Current OCV systems require accurately labelled data to be utilised for training, but the labelling process is time-consuming, expensive and requires expertise. They also rely on consistency in date code format, packaging and camera view angle which is difficult to ensure in a manufacturing environment, so there is a great need for a more robust solution.

## 1.2 Contribution

We propose a deep Bayesian self-training methodology orthogonal to [2] that leverages approximate variational inference in DNNs to estimate predictive uncertainty during a self-training setting. Both aleatoric and epistemic uncertainties of predicted pseudo-labels for unseen data are estimated, and the samples with the lowest predictive uncertainty (highest confidence) are added to the training set in an automated manner. We offer ways to mitigate the known problem of propagating errors in self-training by including: (1) an entropy penalty on the log-likelihood loss to punish overconfident output distributions and facilitate thresholding, and (2) an adaptive sample-wise weight on the influence of predicted pseudo-labelled samples over gradient updates to be inversely proportional to their predictive uncertainty. Lastly, we propose a new simple methodology for visualising and analysing variability between two dataset distributions in DNNs and attempt to

adapt information from one problem to the other by clustering learnt latent variable representations in the context of our application domain. An experimental study on both public and private (real) datasets is presented demonstrating the increased performance of our algorithm over standard self-training baselines.

## 2 Related work

Deep learning model's ability to learn abstract hierarchical representations from data has pushed the state of the art in most machine learning-related tasks [1, 15]. The uptake of these methodologies in academia and industry has resulted in many diverse and interesting DNN applications, wherein patterns learned from data have been adapted to perform tasks in various domains, including computer vision [13, 15–17], medical imaging [18–20] and signal processing [21, 22]. Although many important improvements to DNNs have been made in various domains, there are still many adversities in training models which can be easily adapted to other tasks; and the lack of annotated data is one of the contributing factors.

### 2.1 Deep semi-supervised learning

Most related work addressing the aforementioned issues is often related to domain adaptation philosophy and semi-supervised learning algorithms such as self-training [4], which is an iterative procedure for self-labelling data points in an unlabelled pool, and retraining a classifier until stop conditions are met. Co-training [6] can be considered multiview variant of self-training wherein two separate classifiers are trained on different views of the data and augment each others training sets with their predicted labels. Tri-training [8] extends co-training by having three classifiers, and unlabelled examples are added to a classifier's training set iff the other two agree on the predicted label. Active learning [7] selects the most informative samples from a pool of unlabelled data and retrains the classifier with human given labels in an effort to maximise performance and minimise data labelling requirements. Transfer learning [5] is often used when there is a lack of annotated data in the target domain, and the goal is to adapt knowledge from one task to another by initialising the weights of the target task with the pre-trained weights of another, often performing better than random initialisation. Among these algorithms, transfer learning has undoubtedly had the most success in the context of deep models, and it is widely used in computer vision for adapting visual features from large source domains, to target domains with limited annotated data. Notably, [11] find that initialising a network with transferred features boosts generalisation that

lingers even after fine-tuning to the target dataset, and transferring features from distant tasks is still better than using random weights. Recent work in [23] suggests that a single DL model can jointly learn a number of tasks from multiple domains successfully. In fact, it was observed that adding knowledge from unrelated tasks never hurts performance, rather mostly improves it on all tasks. This phenomenon is complimented by research in [24], with results suggesting that combining tasks, even via a naïve multihead architecture, always improves performance. Authors in [25] propose learning a network comprised of the most successful layers from many different source networks, which are continuously generated and evaluated by a recurrent neural network (RNN) controller. Task transfer learning was recently studied in great depth by [12], where a fully computational approach termed taskonomy was proposed. This was achieved by identifying dependencies between 26 different tasks in latent space, producing a computational taxonomic map for task transfer learning. Deep generative modelling is also gaining popularity in tackling adaptation of knowledge learnt from data generating distributions to pool sets of unlabelled data [26–28]. Other notable related works presented more recently include co-teaching [10], wherein two neural networks are trained simultaneously and teach each other to select clean labels and then decide what data to use for training. Mean teacher models [29] maintain an exponential moving average of model weights and penalise inconsistent predictions, enabling training with fewer labels as an added benefit. Deep co-training [30] extends the original co-training algorithm by training multiple DNNs with different views generated by exploiting adversarial examples. In [31], a simple method termed pseudo-label similar to entropy regularisation [32] is proposed, and it consists of iteratively assigning pseudo-labels via the maximum predicted probability of a NN. Although research on self-training with deep models is scarce, notable work in [33] presents an unsupervised domain adaptation (UDA) framework based on self-training for semantic segmentation using DNNs. They develop a self-paced policy that increases the number of pseudo-labels incorporated in each additional round and demonstrate performance benefits over other popular methods. However, as is the case with all previous works mentioned thus far, their proposed approach does not provide principled predictive uncertainty estimates. The black box nature of DNNs is a concern in most real-world applications, and by quantifying what a model does not know with uncertainty measures, we can not only better trust our predictions but also avoid potentially harmful outcomes [34]. With that in mind, perhaps the most significant related work is in [2], where the authors propose a Bayesian formulation of active learning for image data using DNNs, obtaining a significant

improvement on existing active learning approaches by considering uncertainty estimates in approximating acquisition functions.

## 2.2 Uncertainty estimation

The estimation of uncertainty as a measure of confidence over a model's predictions is desirable for self-labelling, and for safety-critical systems in general [34]. Bayesian neural networks (BNNs) were studied by many in the past [35–37] and have more recently regained popularity. In BNNs, uncertainty is typically captured by placing a prior distribution such as a Gaussian, over the weights and averaging over all possible parameters, rather than optimising them directly. Bayesian inference is then used to compute the posterior over the weights capturing the set of likely parameters. However, BNNs are difficult to perform inference in with traditional methods, as they do not scale well scale to high-dimensional inputs or very complex DL models [34]. Recent promising methods including [34, 38, 39] offer alternative ways of capturing uncertainty by simple modifications to loss functions, having the network to learn/predict aleatoric uncertainty in an unsupervised manner. Aleatoric uncertainty relates to sensory noise in the acquisition process of the data and is therefore inherently irreducible [40]. However, we argue that it can be a great tool for quantifying our uncertainty about pseudo-label predictions. In [39], dropout was shown to perform approximate variational inference, wherein stochastic forward passes with dropout at test time are effectively samples from the approximate posterior. This technique is known as Monte Carlo (MC) dropout [39] and can be used to quantify epistemic uncertainty in NN predictions. Epistemic uncertainty relates to our uncertainty about the model parameters, which is in fact reducible as we observe more data. This is because we can explain the uncertainties about the model parameters in the limit of observing all explanatory variables of the data [34, 40]. This type of uncertainty is useful for identifying out-of-distribution data points and is the most important type of uncertainty measure when assigning pseudo-labels to data.

In this paper we argue that with some modifications, uncertainty estimation techniques in Bayesian deep learning can also be useful in a self-training setting, and to the best of our knowledge, these ideas have yet to be explored in this context. All things considered, we propose a deep Bayesian self-training algorithm, in which a DNN assigns pseudo-labels to new data and automatically weighs their sample-wise importance for the next self-training iteration to be inversely proportional to the predictive uncertainty of the assigned pseudo-label. In this way, we can reduce the burden of manual data annotation requirements and also

offer a measure of uncertainty about our predictions which is important in safety-critical domains.

### 3 Deep Bayesian self-Training

In this section, we provide a brief background on Bayesian NNs and explore the idea of uncertainty estimation of pseudo-label predictions for unlabelled data, in a deep Bayesian self-training framework (see Algorithm 1). In order to quantify what our algorithm does and does not know, we extend existing approaches for estimating uncertainty in deep CNNs [34, 41]. To this end, we consider the following Bayesian formulation of a deep CNN for estimating both aleatoric and epistemic uncertainties.

#### 3.1 Bayesian neural networks

Let  $\mathcal{D} = \{(\mathbf{X}, \mathbf{Y})\}$  denote a dataset given as  $N$  pairs of inputs  $\mathbf{x}_i \in \mathbb{R}^d$  of dimension  $d$ , and class labels  $\mathbf{y}_i \in \{1, \dots, K\}$  of  $K$  total classes. Assuming a Bayesian neural network (BNN) formulation, we place a Gaussian prior probability distribution  $p(\omega)$  over the set of trainable parameters  $\omega = \{\mathbf{W}_1, \dots, \mathbf{W}_\ell\}$ . We define the likelihood conditional output distribution  $p(\mathbf{Y}|\mathbf{X}, \omega)$  of NN for mapping inputs to labels, by finding parameters  $\omega$  that yield the maximum likelihood estimate (MLE). MLE is the pillar of supervised learning in DNNs and is defined as

$$\hat{\omega}_{\text{ML}} = \arg \max_{\omega} \sum_{i=1}^N \log p(\mathbf{y}_i | \mathbf{x}_i, \omega), \quad (1)$$

yielding a point estimate for the most likely parameters to have generated the data. In a Bayesian sense, the MLE is a special case of maximum a posteriori (MAP) estimation when a uniform prior is assumed. In practical classification tasks, the MLE estimator is obtained by minimising the negative log-likelihood of a Bernoulli or softmax distribution depending on the number of classes. We define the softmax negative log-likelihood of our classification NN model as

$$-\log p(\mathbf{y}_i = k | \mathbf{x}, \omega) = -\left(\mathbf{z}_k - \log \sum_{k'} \exp(\mathbf{z}_{k'})\right) \quad (2)$$

where  $\mathbf{z}$  denotes the vector of output logits by the network and  $k$  denotes a class. Having defined a prior and a likelihood, we would like to compute the posterior probability distribution over the weights given the data by Bayes rule

$$p(\omega | \mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y} | \mathbf{X}, \omega) p(\omega)}{p(\mathbf{Y} | \mathbf{X})} \propto p(\mathbf{Y} | \mathbf{X}, \omega) p(\omega), \quad (3)$$

with which we can also formulate the predictive distribution given new inputs  $\mathbf{x}^*$  and labels  $\mathbf{y}^*$

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \omega) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega, \quad (4)$$

enabling predictions using a full distribution over the parameters  $\omega$ , which captures uncertainty over the model parameters, rather than using a point estimate. However, in most cases, the posterior distribution  $p(\omega | \mathbf{X}, \mathbf{Y})$  cannot be evaluated analytically. This is because to compute the marginal probability  $p(\mathbf{Y} | \mathbf{X})$  we must integrate over all possible model parameters  $\omega$  with weighted probability  $p(\omega)$ , in order to obtain the normalising constant, also known as the model evidence. Since the true posterior distribution  $p(\omega | \mathbf{X}, \mathbf{Y})$  is intractable, various approximations exist [36, 37, 42]. Most of them were important early steps towards performing approximate inference in Bayesian NNs, but are unfortunately difficult to employ in modern applications due to scalability constraints or expert knowledge requirements. More recent work in [41, 43–45] addressed some of these issues with variational inference, reigniting interest in the field of Bayesian NNs.

#### 3.2 Variational inference

Next, we provide a background on variational inference (VI) to contextualise some of the ideas presented in [41], wherein dropout is shown to perform approximate variational inference in NNs when used at test time. In VI, a factorised variational distribution from a tractable family  $q_\theta(\omega)$ , parameterised by  $\theta$ , is defined for approximating the posterior distribution by minimising the Kullback–Leibler (KL) divergence between  $q_\theta(\omega)$  and  $p(\omega | \mathbf{X}, \mathbf{Y})$ . Intuitively, the KL divergence is a non-negative asymmetric measure of similarity between the two distributions  $\text{KL}(q_\theta(\omega) || p(\omega | \mathbf{X}, \mathbf{Y}))$ , which we minimise via the variational parameters  $\theta$  of our approximating distribution  $q_\theta(\omega)$

$$\hat{\theta} = \arg \min_{\theta} \mathbb{E}_{q_\theta(\omega)} [\log q_\theta(\omega) - \log p(\omega | \mathbf{X}, \mathbf{Y})]. \quad (5)$$

However, optimising the KL divergence directly requires knowledge of the intractable posterior. This is circumvented by instead maximising the evidence lower bound (ELBO) on the marginal log-likelihood  $\log p(\mathbf{Y} | \mathbf{X})$ , derived via Jensen's inequality  $\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$

$$\mathcal{L}_{\text{ELBO}}(\theta) = \log p(\mathbf{Y} | \mathbf{X}) - \text{KL}(q_\theta(\omega) || p(\omega | \mathbf{X}, \mathbf{Y})), \quad (6)$$

and given that the KL divergence  $\geq 0$  then

$$\log p(\mathbf{Y} | \mathbf{X}) = \mathcal{L}_{\text{ELBO}}(\theta) + \text{KL}(q_\theta(\omega) || p(\omega | \mathbf{X}, \mathbf{Y})). \quad (7)$$

By maximising the lower bound, we implicitly maximise  $\log p(\mathbf{Y} | \mathbf{X})$  and minimise the KL divergence as intended. We extend these ideas in the light of recent developments

in [41] with the Monte Carlo dropout approximation using  $q_\theta(\omega)$ , further explained in the following section.

regularisation term which constrains the approximate posterior  $q_\theta(\omega)$  from deviating too far from prior  $p(\omega)$ . Following [38], we can approximate the KL term with

---

**Algorithm 1** Deep Bayesian Self-Training
 

---

```

1: Note: Pseudo code for training a progressively growing DenseNet in a Bayesian Self-Training setting.
   The incremental growth factor  $\nu$  is adjusted for dataset size.
2: function DBST( $\mathcal{D} = \{(\mathbf{x}, \mathbf{y}, \lambda)\}_{i=1}^N$ ,  $\mathcal{U} = \{\tilde{\mathbf{x}}\}_{i=1}^{\tilde{N}}$ ) ▷ Input training and unlabelled datasets
3:   Initialise :  $r \leftarrow 0$ ,  $k \leftarrow 12$ ,  $k_{\max} \leftarrow 24$ 
4:   Initialise  $\forall \mathbf{x}_i \in \mathcal{D}$  :  $\lambda_i \leftarrow 1$ 
5:   while  $|\mathcal{U}| > 0$  do ▷ Unlabelled dataset cardinality
6:      $r \leftarrow r + 1$ ,  $k \leftarrow \min(k + \nu \cdot (r - 1), k_{\max})$ 
7:      $f(\mathbf{x}) \leftarrow \text{TRAIN}(\mathcal{D}, k)$  ▷ Train a DenseNet with growth rate  $k$ 
8:      $(\hat{\mathbf{p}}, \hat{\mathbf{s}}) \leftarrow \text{MC DROPOUT}(f(\mathbf{x}), \mathcal{U})$ 
9:     for  $\tilde{\mathbf{x}}_i \in \mathcal{U}$  do
10:       $\text{Var}[\mathbf{y}_i] \leftarrow \exp(\hat{\mathbf{s}}_i) + \mathbb{H}[\hat{\mathbf{p}}_i]$  ▷ Aleatoric and Epistemic uncertainties
11:       $\hat{\mathbf{y}}_i \leftarrow \arg \max \hat{\mathbf{p}}_i$ 
12:      if  $\text{Var}[\mathbf{y}_i] < \tau$  then ▷  $\tau$  is computed via IQR
13:         $\lambda_i \leftarrow 1 / \exp(\text{Var}[\mathbf{y}_i])^{\phi(\tau)}$ 
14:         $\mathcal{D} \leftarrow \mathcal{D} \cup \{\tilde{\mathbf{x}}_i, \hat{\mathbf{y}}_i, \lambda_i\}$  ▷ Add weighted pseudo-labelled sample to  $\mathcal{D}$ 

```

---

### 3.3 Continuous relaxation of dropout

Concrete dropout is based on concrete relaxation of discrete distributions [46], allowing the replacement of dropout's discrete Bernoulli distribution with its continuous relaxation [47]. To obtain calibrated uncertainty estimates with Monte Carlo dropout, it is necessary to tune the dropout probabilities. A grid search is a common but costly approach for large models, highlighting the benefit of optimising them directly with gradient descent. This requires formulating an objective for minimising epistemic uncertainty [41] using the variational interpretation of dropout.

Formally, dropout can be treated as an approximating distribution  $q_\theta(\omega)$  to the posterior in a BNN, where  $\omega$  represents the weight matrices of the  $\ell$ th of  $L$  layers in the network  $\omega = \{\mathbf{W}_\ell\}_{\ell=1}^L$ , and  $\theta$  are the variational parameters to optimise [47]. Let  $\mathcal{F}(\omega)$  be the model with weight matrix realisation  $\omega$ ; given a random set  $S$  comprising  $M$  of all  $N$  data points, denote the model's output on the  $\mathbf{x}_i$  input as  $\mathcal{F}(\mathbf{x}_i; \omega)$ . The following NN objective function can then be formulated

$$\begin{aligned} \mathcal{L}_{\text{MC}}(\theta) = & -\frac{1}{M} \sum_{i \in S} \log p(\mathbf{y}_i | \mathcal{F}(\mathbf{x}_i; \omega)) \\ & + \frac{1}{N} \text{KL}(q_\theta(\omega) \parallel p(\omega)), \end{aligned} \quad (8)$$

where  $p(\mathbf{y}_i | \mathcal{F}(\mathbf{x}_i; \omega))$  is the model's likelihood, a Gaussian with a predictive mean given by  $\mathcal{F}(\mathbf{x}_i; \omega)$ . KL is a

$$\text{KL}(q_M(\mathbf{W}) \parallel p(\mathbf{W})) \propto \frac{I^2(1-p)}{2} \|\mathbf{M}\|^2 - K\mathbb{H}[p], \quad (9)$$

where  $\{\mathbf{M}_\ell, p_\ell\}_{\ell=1}^L$  is a set of mean weight matrices and dropout probabilities, such that (s.t.)  $q_{M_\ell}(\mathbf{W}_\ell) = \mathbf{M}_\ell \cdot \text{diag}[\text{Bernoulli}(1-p_\ell)^{K_\ell}]$  for a single NN weight matrix  $\mathbf{W}_\ell \in \mathbb{R}^{K_{\ell+1} \times K_\ell}$ .  $\mathbb{H}[p]$  is simply the entropy of a Bernoulli random variable with probability  $p$

$$\mathbb{H}[p] := -p \log p - (1-p) \log(1-p), \quad (10)$$

which can be interpreted as a regularisation term that only depends on dropout probability  $p$ , so minimising the KL term is equivalent to maximising the entropy of a Bernoulli random variable with probability  $(1-p)$ . Rather than sampling the random variable from the discrete Bernoulli distribution, by adopting the concrete distribution [46, 47] with some temperature  $t$ , it is possible to sample variables in the interval  $[0, 1]$ , s.t. the concrete relaxation distribution  $\tilde{\mathbf{z}}$

$$\begin{aligned} \tilde{\mathbf{z}} = \text{sigmoid} \left( \frac{1}{t} \cdot [\log p - \log(1-p) \right. \\ \left. + \log u - \log(1-u)] \right), \end{aligned} \quad (11)$$

parameterised by means of  $u \sim \text{Unif}(0, 1)$ , provides a relationship between  $\tilde{\mathbf{z}}$  and  $u$ , which is differentiable w.r.t.  $p$ . With the concrete relaxation of the dropout masks, the dropout probabilities for each layer  $\{p_\ell\}_{\ell=1}^L$  can be optimised using the path-wise derivative estimator [47].



### 3.4 Entropy penalty on output distributions

The probabilities assigned to incorrect classes at test time help quantify a model's ability to generalise. By penalising output distributions with low entropy (i.e. confident predictions), we can obtain a similar effect to label smoothing and improve generalisation [48]. This can be useful in self-training, since we assign pseudo-labels based on low uncertainty predictions, which are in some cases wrongly assigned. We suggest that by penalising very confident output distributions we can improve generalisation and make thresholding easier since the output distributions are smoother, rather than overly concentrated at 0 or 1. The entropy of a NNs output conditional distribution is given by

$$\mathbb{H}[p(\mathbf{y}|\mathbf{x}, \omega)] = - \sum_i p(\mathbf{y}_i|\mathbf{x}, \omega) \log p(\mathbf{y}_i|\mathbf{x}, \omega), \quad (12)$$

with  $p(\mathbf{y}|\mathbf{x}, \omega)$  as the probability distribution obtained from a softmax function. To penalise very confident predictions, we can simply take the negative log-likelihood and subtract the entropy of the output distribution as

$$\mathcal{L}_{\text{NLL}}(\omega) = - \sum \log p(\mathbf{y}|\mathbf{x}, \omega) - \beta \mathbb{H}[p(\mathbf{y}|\mathbf{x}, \omega)], \quad (13)$$

where the scaling hyperparameter  $\beta$  balances how much we would like to penalise non-uniformity of the softmax.

### 3.5 Inverse uncertainty weighting

A known limitation of self-training is the potential accumulation of wrongly pseudo-labelled samples being added to the training set. A common approach is to remove less confident samples from the training set and leave them in the unlabelled set. However, this tends to underperform in practice, as the algorithm can become biased by continuously adding the easiest unlabelled samples to the training set. This can hinder learning over time, as more difficult and potentially informative samples are neglected.

In attempt to mitigate this behaviour, we propose a sample-wise weighting scheme during training that places a weight on each training sample  $\{\mathbf{x}_i, \hat{\mathbf{y}}_i, \lambda_i\}$ , proportional to the predictive uncertainty over its pseudo-label  $\hat{\mathbf{y}}_i$ , such that its contribution to the loss function is inversely proportional to its predictive uncertainty (see Algorithm 1). To calculate the predictive uncertainty, we can have the network predict the aleatoric uncertainty as one of its outputs and add the epistemic uncertainty obtained from the variance of Monte Carlo dropout samples.

Formally, let  $\hat{\mathbf{p}}_t = \text{softmax}(\mathcal{F}(\mathbf{x}; \hat{\omega}_t))$  denote the softmax out of a BNN, and  $\{\hat{\mathbf{p}}\}_{t=1}^T$  be the set of outputs from  $T$  Monte Carlo dropout samples at test time, each parameterised by weights drawn from the approximate posterior

$\hat{\omega}_t \sim q_{\hat{\theta}}(\omega)$ . We propose calculating the predictive uncertainty from these samples by generalising the binary variant approach in [49] to a multivariate classification setting. By the definition of variance of a multinomial distribution, we can decompose the variance of  $\hat{\mathbf{p}}$  into

$$\text{Var}[\hat{\mathbf{p}}] \approx \text{tr} \left( \mathbb{E}[\text{diag}(\hat{\mathbf{p}}) - \hat{\mathbf{p}}\hat{\mathbf{p}}^\top] + \mathbb{E}[\hat{\mathbf{p}}^2] - \mathbb{E}[\hat{\mathbf{p}}]^2 \right), \quad (14)$$

where the first term represents aleatoric uncertainty  $\sigma_a^2$ , and the second is the epistemic  $\sigma_e^2$ . Each diagonal entry of the resulting matrix is the variance of a binomially distributed random variable, and the off-diagonals are negative covariances for fixed  $T$ . Since we are only interested in a single number to measure our uncertainty, we take trace of the resulting uncertainty matrix.

Alternatively, we can have the NN predict the input noise variance  $\sigma_a^2$  as one of its outputs [34], by assuming measurement error in our target function  $y = \mathcal{F}(\mathbf{x}) + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \sigma_a^2)$ . The predictive variance in a multivariate classification setting is then given by

$$\begin{aligned} \text{Var}[\hat{\mathbf{y}}] &\approx \frac{1}{T} \sum_t \exp(\hat{\mathbf{s}}_t) - \sum_j \mathbb{E}[\hat{\mathbf{p}}] \log \mathbb{E}[\hat{\mathbf{p}}] \\ &= \mathbb{E}[\exp(\hat{\mathbf{s}})] + \mathbb{H}[\text{softmax}(\mathcal{F}(\mathbf{x}; \hat{\omega}))], \end{aligned} \quad (15)$$

the entropy term measures epistemic uncertainty in the output softmax distributions, whereas the log aleatoric uncertainty  $\hat{\mathbf{s}}_i := \log \sigma_{a,i}^2$  term is regressed by the NN for each input  $\mathbf{x}_i$ , for numerical stability. To capture aleatoric uncertainty in our classification task, we can use Monte Carlo integration on the NNs Gaussian log-likelihood objective function, by drawing  $t \in T$  samples of Gaussian noise-corrupted NN output logits  $\mathcal{F}(\mathbf{x})$ , yielding the following loss

$$\mathcal{L}_{\text{NLL}} = - \log \mathbb{E}[\text{softmax}(\mathcal{F}(\mathbf{x}) + \epsilon_t \odot \exp(\hat{\mathbf{s}}|\mathbf{x}))], \quad (16)$$

with  $\epsilon_t \sim \mathcal{N}(0, I)$  parameterised by the predicted aleatoric uncertainty  $\exp(\hat{\mathbf{s}}|\mathbf{x})$  for each sample  $\mathbf{x}_i$ , which learns to capture measurement error.

Having calculated the predictive uncertainty  $\text{Var}[\hat{\mathbf{p}}]$  of our pseudo-labels, we calculate a per-sample importance weight  $\{\mathbf{x}_i, \hat{\mathbf{y}}_i, \lambda_i\}$  with

$$\lambda_i = \frac{1}{\exp(\text{Var}[\hat{\mathbf{p}}_i])^{\phi(r)}}, \quad (17)$$

where  $\phi(\cdot)$  is a parameterised hyperbolic tangent function

$$\phi(r) = \frac{1 - \exp(\gamma \cdot r + b)}{1 + \exp(\gamma \cdot r + b)}, \quad (18)$$

with  $\gamma, b$  as scale and intercept terms, and  $r$  denotes the self-training iteration. The weighted penalised log-likelihood of our NN with weights  $\omega$  is then

$$\mathcal{L}_{\text{PLL}} = \sum_i \lambda_i \log p(y_i | \mathbf{x}_i, \omega) - \beta \mathbb{H}[p(y_i | \mathbf{x}_i, \omega)], \quad (19)$$

where  $p(y|\mathbf{x}, \omega)$  is computed via softmax, and the optional confidence entropy penalty term is balanced by  $\beta$ . By tuning  $\gamma$  and  $b$ , we can obtain the desired behaviour over  $r$  iterations, s.t. when the uncertainty is low, we assign high weight to the predicted pseudo-labelled sample  $\{\mathbf{x}_i, \hat{y}_i, \lambda_i \approx 1\}$ . We can incrementally encourage the model to assign more weight to uncertain pseudo-labelled samples as self-training progresses, since in the  $\lim_{r \rightarrow \infty} \phi(r) = -1$ . Intuitively, this procedure inverts Eq. (17) over time, incrementally forcing exploration by adding more uncertain, and potentially informative samples, to the training set. In summary, using this logic along with entropy penalties on overconfident output distributions, we can mitigate the effect of pseudo-labelling error accumulation in the training set and adjust risk taking by tuning  $\gamma$  and  $b$ . Once per-sample predictive uncertainties are calculated, we decide on which pseudo-labelled samples to add to the training set via a Tukey fence. Intuitively, assume a NN has been trained on data  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , learning a function  $\mathcal{F}(\mathbf{x}; \omega)$  for mapping inputs to labels. At inference time, we take the correct predictions where  $y_i = \mathcal{F}(\mathbf{x}_i; \omega)$  and retrieve their predictive uncertainty. We then summarise variability by calculating the interquartile range (IQR) outlier statistic and define an uncertainty upper bound  $\tau$ , which is used to decide which pseudo-labelled samples from  $\mathcal{U} = \{\tilde{\mathbf{x}}_i\}_{i=1}^{\tilde{N}}$  should be added to  $\mathcal{D}$  following

$$\mathcal{D}^* = \forall_i \in \{\mathcal{D} \cup \{\tilde{\mathbf{x}}_i, \hat{y}_i, \lambda_i\} \mid \text{Var}[p(y_i | \mathbf{x}_i)] < \tau\}, \quad (20)$$

where  $\hat{y}_i$  denotes the pseudo-label assigned to sample  $\mathbf{x}_i$  computed as  $\hat{y}_i = \arg \max \hat{\mathbf{p}}_i$ , and  $\mathcal{D}^*$  is the augmented training set. Lastly, we can also easily adjust the uncertainty upper bound  $\tau$  by selecting higher or lower quartiles to reflect how confident we would like to be about predictions before adding samples to  $\mathcal{D}^*$ .

## 4 Latent variable adaptive clustering

We propose a new simple methodology for visualising and analysing variability between distributions and attempt to adapt information from one problem to another in DNNs. In Fig. 1, an illustration of our adaptation framework is shown using an example backbone InceptionV3 CNN. Let the following denote two training sets from separate datasets targeting the same task

$$\begin{aligned} \mathcal{D}_1 &= \{(\mathbf{x}_1^{(i)}, \mathbf{y}_1^{(i)}) ; i = 1, \dots, N_1\}, \\ \mathcal{D}_2 &= \{(\mathbf{x}_2^{(i)}, \mathbf{y}_2^{(i)}) ; i = 1, \dots, N_2\}, \end{aligned} \quad (21)$$

and the two respective test sets as

$$\begin{aligned} \mathcal{T}_1 &= \{(\tilde{\mathbf{x}}_1^{(i)}, \tilde{\mathbf{y}}_1^{(i)}) ; i = 1, \dots, \tilde{N}_1\}, \\ \mathcal{T}_2 &= \{(\tilde{\mathbf{x}}_2^{(i)}, \tilde{\mathbf{y}}_2^{(i)}) ; i = 1, \dots, \tilde{N}_2\}. \end{aligned} \quad (22)$$

Let  $\mathcal{F}(\mathcal{D}_1; \mathbf{W}_1)$  and  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$  denote two architecturally identical CNNs trained separately on each dataset. For each CNN, we extract the final fully connected layer activations  $\{\mathbf{x}_1^{(i)}, \tilde{\mathbf{x}}_1^{(i)}\} \in \mathbb{R}^{2048}$  and  $\{\mathbf{x}_2^{(i)}, \tilde{\mathbf{x}}_2^{(i)}\} \in \mathbb{R}^{2048}$  as latent variables representations, by simply forward-propagating each image through as is typically done at inference time.

Utilising these, our adaptation methodology is then performed as follows:

1. Given  $\mathcal{D}_2$ , produce a set of clusters  $\mathbf{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_k\}$  by minimising the within-cluster  $L^2$  norms of the following clustering objective function

$$\hat{\mathbf{C}}_{k\text{-means}} = \arg \min_{\mathbf{C}} \sum_{i=1}^k \sum_{\mathbf{x} \in \mathbf{C}_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2. \quad (23)$$

2. Repeat step 1 with  $\mathcal{D}_1$  to generate  $k$  clusters  $\mathbf{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  and compute the  $k$  closest instances in  $\mathcal{D}_1$  to each centroid in  $\mathbf{U}$ . Fetch the corresponding set of images  $\mathbf{S} = \{\mathbf{S}_1, \dots, \mathbf{S}_k\}$ , whose latent variables are closest to  $\mathbf{U}$ ;
3. Forward-propagate  $\mathbf{S}$  through  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$  to obtain a new set of adapted clusters  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_k\}$ , where  $\mathbf{S}$  is considered an approximation of  $\mathbf{U}$  from  $\mathcal{F}(\mathcal{D}_1; \mathbf{W}_1)$ ;
4. Derive an augmented cluster representation that encapsulates knowledge from both facets of the trained CNNs, by concatenating the respective  $\mathbf{C}$  and  $\mathbf{Z}$  clusters into a set  $\mathbf{A} = \{\mathbf{c}_1, \dots, \mathbf{c}_k, \mathbf{z}_1, \dots, \mathbf{z}_k\}$ ;
5. Compute the Euclidean distance between  $\mathcal{T}_1$  and  $\mathbf{A}$  and evaluate the classification performance;
6. Iteratively remove the lowest performing cluster in  $\mathbf{A}$  and repeat step 5 until the performance stops improving.

In all cases, the  $k$ -means++ [50] seeding strategy was used, whereby the first cluster centre  $\mathbf{c}_1$  is chosen uniformly at random from  $\mathcal{X}$ , and all preceding cluster centres  $\mathbf{x} \in \mathcal{X}$  are chosen with probability

$$\mathbf{c}_i = \frac{D(\mathbf{x})^2}{\sum_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x})^2}, \quad (24)$$

where  $D(\mathbf{x})$  denotes the distance between  $\mathbf{x}$  and the closest  $\mathbf{c}_i$ . Moreover, we assign the class label of a given cluster  $\mathbf{c}_i$  as simply the mode class  $j$  of all data points within it

$$\mathbf{c}_i^j = \max_{j \in J} |\mathbf{c}_i \cap j|. \quad (25)$$

In the experimental study of Sect. 6, we demonstrate that our method distils and adapts knowledge from both trained

CNNs on real data, achieving better performance than direct inference of  $\mathcal{T}_1$  with  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$ , without any parameter retraining.

## 5 Experimental study

This section is divided into two separate subsections: the first subsection presents experiments using deep Bayesian self-training applied to the MNIST public domain dataset. An ablation study is presented and comparisons are made with baseline methods. The second subsection comprises a study using private (real) datasets, in which we perform some preliminary experiments using transfer learning and then we evaluate our proposed latent variable adaptable clustering method. We then finish off the second subsection by evaluating deep Bayesian self-training on the self-annotation of the real datasets.

### 5.1 MNIST dataset

In order to validate our algorithm, we conduct a series of self-labelling experiments on the popular MNIST dataset. The MNIST dataset is comprised of 60,000 training and 10,000 testing handwritten digit examples, respectively. Firstly, we try to create a realistic scenario by splitting the 60,000 training examples into a smaller but balanced training set of only 50 examples per class and a validation set of 500 training examples per class and allocate all remaining data to the unlabelled pool set. We begin by defining our backbone NN architecture of choice as a DenseNet [15]. DenseNets have revealed several well-founded advantages over previous architectures, from mitigating vanishing gradients to encouraging feature propagation and reuse with shorter connections between layers [15, 51]. The dense connectivity in DenseNets can be formally defined as

$$\mathbf{A}^{[\ell]} = f\left(\text{BN}\left(\mathbf{W}^{[\ell]} \cdot [\mathbf{A}^{[0]}, \mathbf{A}^{[1]}, \dots, \mathbf{A}^{[\ell-1]}]\right)\right), \quad (26)$$

where  $f(\cdot)$  is the ReLU activation function,  $\text{BN}(\cdot)$  is batch normalisation [52] and  $[\mathbf{A}^{[0]}, \mathbf{A}^{[1]}, \dots, \mathbf{A}^{[\ell-1]}]$  represents feature map-wise concatenation of all layers preceding  $\ell$ . A sequential composite function consisting of BN, ReLU and  $3 \times 3$  convolution can then be defined as  $H^{[\ell]}$ . Each function  $H^{[\ell]}$  produces  $\omega$  feature maps, known as the growth rate of the network, and each layer  $\ell$  takes as input  $f + \omega \times (\ell - 1)$  total feature maps, where  $f$  denotes the number of channels in the visible layer. To reduce spatial dimensionality of feature maps, a transition layer is introduced between densely connected DenseBlocks. Transition layers in [15] are composed of BN followed by  $1 \times 1$  convolution

and  $2 \times 2$  average pooling with a feature map compression factor  $\theta = 0.5$ .

Following Algorithm 1 closely, we propose a progressively growing NN scheme by starting off with a 40 layer deep DenseNet with a growth rate  $k = 12$ , and incrementally increasing the growth rate (width) of the network as more data are added to the training set. In the first iteration, the network has only 181k parameters to avoid overfitting on the small initial training set, but complexity of the network is incrementally increased in an automated way. As described in greater detail in Sect. 3.5, we employ Monte Carlo dropout at test time to calculate the predictive uncertainty of the assigned pseudo-labels samples. In all cases, we take  $T = 30$  samples, equating to 30 different dropout masks. We compare the performance of our proposed approach with a baseline ensemble method (DEST) similar to [53] for estimating predictive uncertainty, and the vanilla self-training methodology, albeit in a deep learning model, considering only the output probability of the NN as a measure of confidence, similarly to [31]. We also evaluate the effect of our inverse uncertainty weighting scheme, as well as the entropy penalty on confident output distributions on the performance of our Bayesian self-training algorithm.

#### 5.1.1 Training details

In all MNIST experiments, we use the same DenseNet model and hyperparameters for fair comparisons. Specifically, we train the networks using stochastic gradient descent (SGD) with a Nesterov momentum of 0.9, a batch size of 32 and an initial learning rate of 0.1. We train all models for 75 epochs and reduce the learning rate by a factor of 10 at 50 and 75% of the way through training. All models are trained using the same train/valid/test/unlabelled splits, no data augmentation is used aside from simple image standardisation (mean 0 sd. 1), and we take  $T = 30$  Monte Carlo dropout samples to at test time as explained in Sect. 3.5. With regard to the ensemble, we train  $M = 5$  models each initialised with random weights and capture the predictive uncertainty following Eq. (15), but without using dropout at test time. Lastly, the stop conditions can be adjusted depending on the application at hand, but here they were kept consistent in all experiments for fairness of comparison. Specifically, we stipulate that if less than the current batch size number of images are selected to be added to the training set in the next self-training iteration, the algorithm stops.

#### 5.1.2 Ablation study

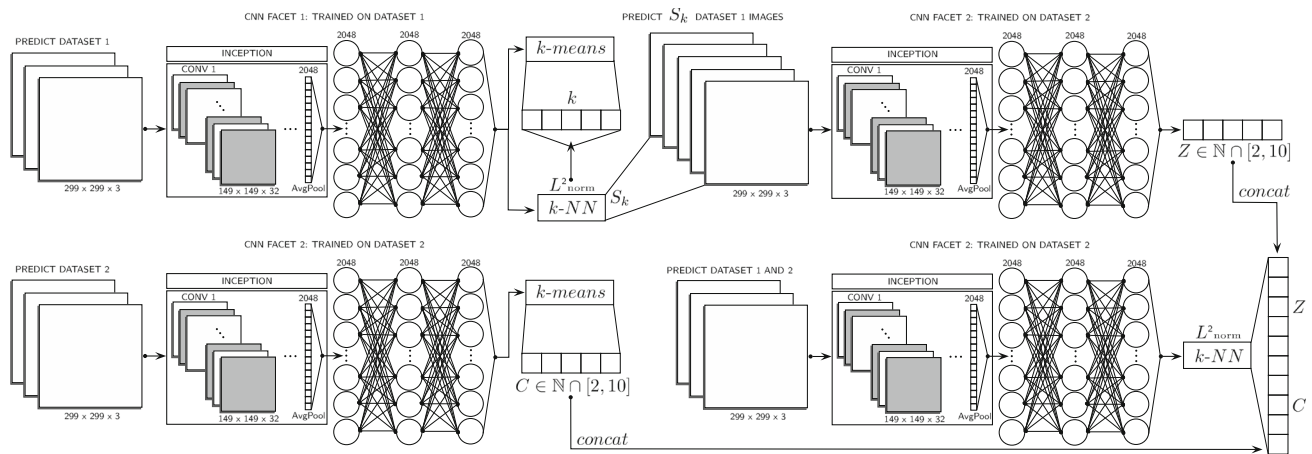
The results are reported in Table 1 and illustrated in Figs. 2, 3, 4 and 5. In our experiments, we simply have the



**Table 1** Deep Bayesian self-training results on self-labelling the MNIST dataset

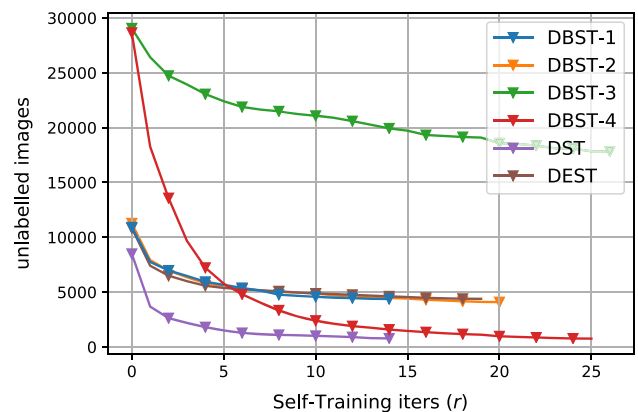
Model	$\tau$	$\lambda_i$	$\mathcal{L}_{\text{PNLL}}$	Precision	Recall	F1-score	Unlabelled	Cohen's $\kappa$	$r$ iters
DST	—	✗	✗	.0103	.0103	.0103	781	.0115	15
DEST	Q3	✓	✗	.0044	.0044	.0044	4391	.0049	20
DBST-1	Q3	✗	✗	.0042	.0043	.0043	5044	.0045	15
DBST-2	Q3	✓	✗	.0032	.0032	.0032	4092	.0035	21
DBST-3	Q2	✓	✗	.001	.001	.001	17,828	.0011	27
DBST-4	Q2	✓	✓	.0071	.007	.0071	762	.0079	26

$\tau$  is the upper bound uncertainty threshold for augmenting  $\mathcal{D}^*$ ,  $\lambda_i$  are sample-wise inverse uncertainty weights, and  $r$  is the number of self-training iterations taken before stop conditions were met. All metrics (precision, recall, F1-score and Cohen's  $\kappa$ ) are reported in 1–metric format



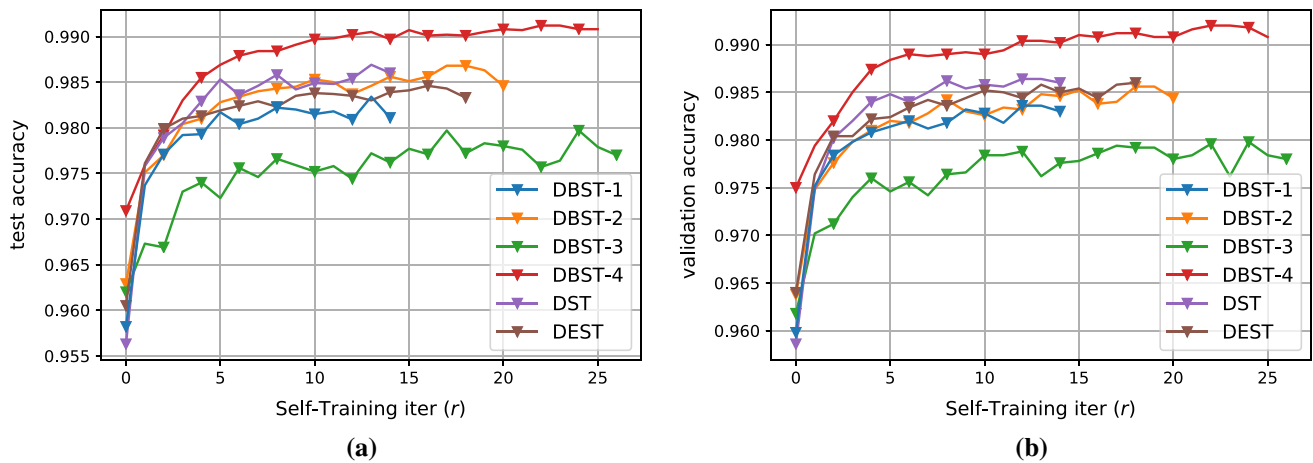
**Fig. 1** Illustration of the multiple CNN facet adaptation framework proposed, which is based on clustering of extracted latent variable representations. The architectural details of each CNN are as described in Fig. 8

NNs predict the labels for the 54,500 unlabelled MNIST samples and evaluate how well the system is doing at predicting the correct labels at the end of each self-training iteration. The evaluation is primarily considered in terms of the Cohen's kappa statistic ( $\kappa$ ) as it is more robust than accuracy by taking into account random luck, and the number of images left unlabelled after self-training. As can be observed from the results, the addition of our proposed inverse uncertainty weighting scheme improves the performance of the algorithm by leaving less images unlabelled and achieving a higher  $\kappa$  score (DBST-1 to DBST-2). We also test the effect of the quartile uncertainty thresholds for  $\tau$  from Q3 to Q2 (DBST-2 to DBST-3), meaning we are more strict about which pseudo-labelled samples we can add to the training set. This only considers very highly confident pseudo-label predictions resulting in a higher  $\kappa$  score, at the cost of labelling less examples as expected. In the DBST-4 model, we combine both the sample-wise inverse uncertainty weighting scheme and the entropy penalty on the log-likelihood loss ( $\mathcal{L}_{\text{PNLL}}$ ) using  $\beta = 1$  as described in Sect. 3.5. As reported in Table 1, the



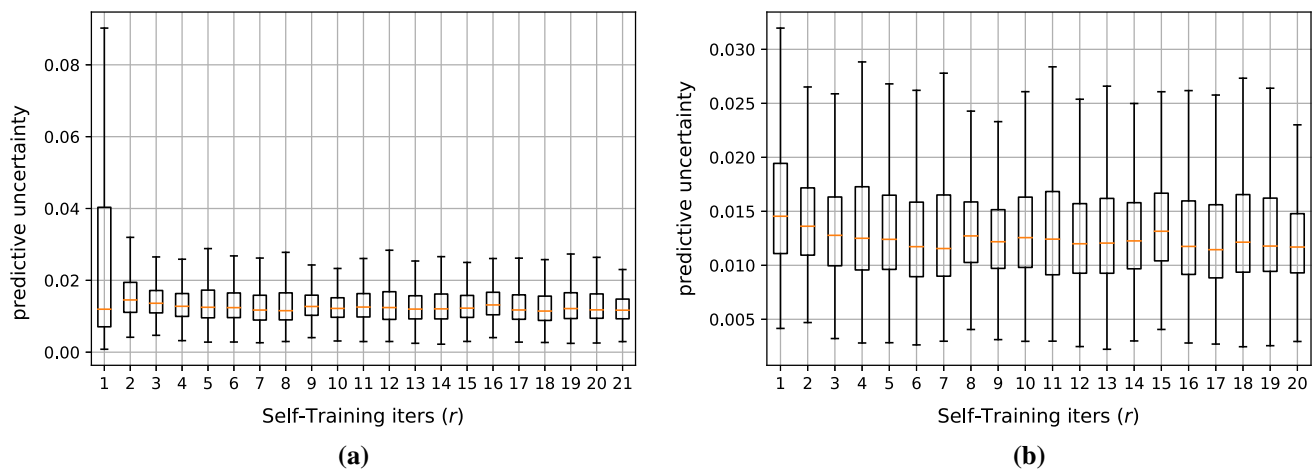
**Fig. 2** Self-training model comparisons regarding number of images left unlabelled after  $r$  iterations. Notice how the baseline self-training (DST) is overconfident by wrongly pseudo-labelling more samples early and propagating these errors, resulting in a lower Cohen's  $\kappa$  score as reported in Table 1

number of examples left unlabelled is significantly less, whilst maintaining a good Cohen's  $\kappa$  agreement between predicted and actual labels. In comparison with the others,



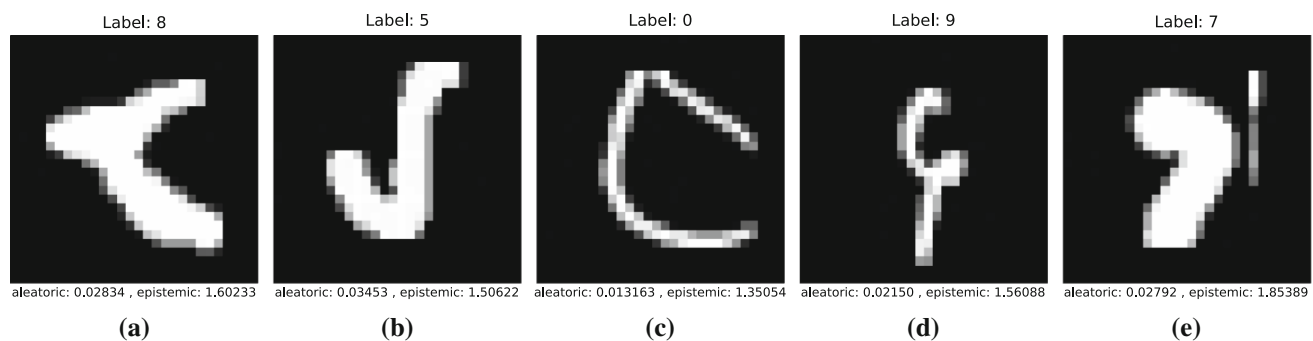
**Fig. 3** Model performance comparisons over  $r$  self-training iterations. **a** MNIST test set performance after each self-training iteration. **b** As in **a** but comparing validation set performance. Notice that every model uses the same stop condition for fair comparison, but they stop

at different times due to their uncertainty level. DBST-4 using both inverse uncertainty sample weights and an entropy penalty on the log-likelihood loss ( $\mathcal{L}_{\text{PNLL}}$ ) generalises better as reported in Table 1



**Fig. 4** Box plots (IQR) depicting the quartiles for setting the uncertainty upper bound threshold  $\tau$  over  $r$  iterations in the DBST-2 model as an example. Note: these IQR stats are calculated using the

predictive uncertainties of correctly classified samples in the train/valid/test sets only. **a** Shows all iterations ( $r = 21$ ) whereas **b** omits the first one for better visibility



**Fig. 5** Examples of images left in the unlabelled pool set for model DBST-2. Images with the highest epistemic uncertainty were selected for each digit class, along with their corresponding aleatoric uncertainties reported in the x-axis. The actual label of each image

is found on top. As we can see from these difficult examples, these digits were automatically identified as problematic (too uncertain) in the DBST pseudo-labelling process, so they were not added to the training set  $\mathcal{D}^*$

the DBST-4 model provides the best balance between the number of unlabelled images left after self-training and a high Cohen's  $\kappa$  score.

### 5.1.3 Comparative discussion

Lastly, we compare our Bayesian models (DBST) with two baseline method for estimating uncertainty in a similar way to [53], known as a deep ensemble of NNs (DEST), and the standard self-training (DST) following the logic in [31], and simply using the NNs predicted probability of an assigned pseudo-label as a level of confidence. The predictions from each NN in the ensemble (DEST) can be used as to calculate predictive uncertainty as the deviations capture model parameter uncertainty. Here, we do not employ any bootstrap methods as the randomness from the NN weight initialisation and shuffled training has been shown to be sufficient experimentally [53]. We use the same DenseNet architecture, including related hyperparameters and identical dataset splits to train an ensemble of five models. Table 1 shows that our methods (DBST) are better than using an ensemble ( $M = 5$ ) for predicting uncertainty for our self-training purpose, whilst taking approximately  $5\times$  less time to run in our experiments. Note that Monte Carlo dropout samples are very cheap to compute at inference time compared to training multiple models; thus, we can afford to take multiple samples, i.e.  $T = 30$  as compared to an ensemble of  $M = 5$ , which is also an advantage of our approach.

With regard to the vanilla self-training baseline (DST), again we use the exact same DenseNet architecture and related hyperparameters for fair comparisons. As previously outlined, in standard self-training we take the NNs predicted probability as a measure of confidence, and to demonstrate the inadequacy of this method, we threshold with a very high confidence probability of .99. This simply means that only pseudo-label predictions above the .99 probability (confidence) threshold in a 10-way softmax (MNIST digit classes) are added to the training set. As reported in Table 1 and Fig. 2, DST underperforms compared to our methods since it is overconfident early on, resulting in the addition of more wrong pseudo-labels to the training set, thus propagating the errors forward. Although the number of images left unlabelled is low, the Cohen's  $\kappa$  score is significantly lower

## 5.2 Real datasets

Four datasets of food package photographs were collected by a leading food company and provided to us for research purposes. The four sets include 1404, 6739, 1154 and 13948 captured images, respectively. In order to produce trainable datasets, a portion of the images was first

**Table 2** Number of images per category in each dataset

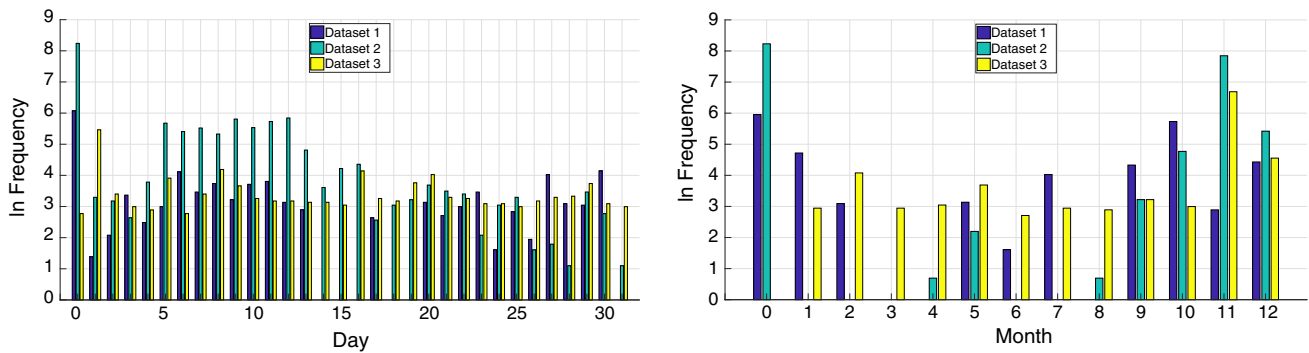
Annotation (DD/MM)	Dataset		
	1	2	3
Missing/missing	375	3715	0
Missing/complete	59	68	16
Complete/missing	10	39	0
Complete/complete	645	2847	1138
Unreadable	315	46	0

manually annotated w.r.t. the presence of *use-by* dates, and lack thereof. In the case of unreadable images, in which dates were not discernible from the background—potentially due to heavy distortion—non-homogeneous illumination or blur was then set aside in a separate category. Conversely, images in which either day or month, or both were missing, were considered as incomplete and subsequently grouped into their own category. Lastly, images of good quality, reporting the date including both the day and month, were considered as good candidates for OCV.

The first three sets of images were annotated as mentioned above to form five categories: complete dates, missing day, missing month, no date and unreadable (Table 2), whereas photographs belonging to the fourth dataset were annotated as good or bad candidates for OCV and utilised to test our proposed Bayesian self-annotating framework. After annotating all the images in the first three datasets, it was possible to plot some statistics (see Fig. 6) on the frequency of specific dates within each dataset, and thus devise a methodology for conducting experiments with balanced sets of classes. Moreover, by inspecting the images with partially missing data, it was observed that most of them were photographs of package labels which had been folded at crucial points, included photographic glare, digits fainting over time, or included human made occlusions. With regard to the fourth dataset, 8931 images were annotated as including readable dates, and the remaining 5017 as unreadable (Fig. 7).

### 5.2.1 Transfer learning

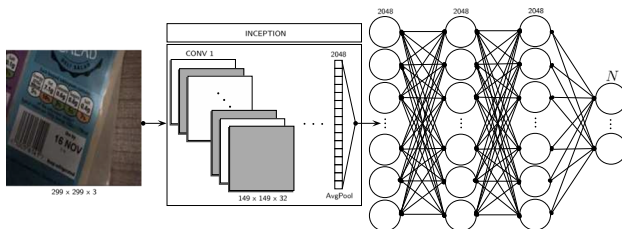
It was of particular interest to conduct transfer learning in order to assess the adaptability of pre-trained CNN weights [54] on the current food datasets. Specifically, each image from our datasets was fed through a previously trained InceptionV3 CNN on the ImageNet dataset, up to the last global average pooling (GAP) layer, where a 2048-dimensional vector representation of each instance was extracted. The 2048-dimensional vectors then became the input to a new series of FC layers and a final softmax layer



**Fig. 6** Left: Frequency (In scale) of appearance per 'Day' in *use-by* dates. Right: Respective appearance per 'Month'



**Fig. 7** Per category examples of images in our datasets. **a** Complete Date (day and month visible). **b** Partial Date (no day visible). **c** Partial Date (no month visible). **d** Unreadable. **e** No date (neither day or month visible)



**Fig. 8** Depiction of the classification architecture. From left to right, input images were resized to  $299 \times 299 \times 3$  to accommodate the CNN's convolutional layer parameters and arithmetic. There exist two hidden layers with 2048 units each and ReLu activations. The number of units  $N$  in the softmax layer was adjusted as per the number of classes being classified in different experiments

able to predict  $N$  classes (see Fig. 8). In order to optimise the training performance of the new FC layer network, a series of architectural decisions were made empirically, and the best performances were achieved using a FC network consisting of two 2048 unit hidden layers with rectified linear unit (ReLU) activations and batch normalisation (BN) [52] layers.

The risk of overfitting rises as the number of parameters increases w.r.t. number of training examples. Due to the limited amount of training data, available for experimentation, it is infeasible to train state-of-the-art models from scratch. Therefore, we introduced an effective regulariser in the new network as well as adapted previously learned low-level features through transfer learning. One of the most effective regularisation techniques is dropout [55]. In

practice, to preserve more information in the input layer  $\ell^{(0)}$  (of  $L$  total layers) in the network and thus aid learning, the probability of keeping ( $p(z^{(i)}) \neq 0$ ) any given neuron  $z^{(i)}$  in layer  $i$  was as defined per the following schema

$$\ell^{(i)} = \begin{cases} p(z^{(i)}) = 0.8 & \text{if } i = 0 \\ p(z^{(i)}) = 0.5 & \text{otherwise.} \end{cases} \quad (27)$$

In view of the unbalance present among the various classes, it was beneficial to use a weighted negative log-likelihood as a loss function (28). In (28),  $\lambda_j$  is a weight coefficient computed for the  $j$ th of all classes  $J$  as a function of the proportion of instances  $N_j$  compared to the most densely populated class (29). During training,  $\lambda$  encourages the model to focus on under-represented classes

$$\mathcal{L}_{\text{NLL}} = - \sum_i \lambda_j y_i \log(\hat{y}_i) - (1 - y_i) \log(1 - \hat{y}_i) \quad (28)$$

calculating the per-class weight parameter  $\lambda_j$  with

$$\lambda_j = \frac{1}{N_j} \max \left( \{N_i\}_{i=[1:J]} \right). \quad (29)$$

In the case of multiclass classification, where  $J > 2$ , the weighted cross entropy loss function can be defined as

$$\mathcal{L}_{\text{NLL}} = - \sum_{i=1}^M \sum_{j=1}^J \lambda_j y_{ij} \log(\hat{y}_{ij}), \quad (30)$$

where  $\log p(\hat{y} = j | z_j)$  is calculated as

**Table 3** Experiment results of OCV binary classification

CNN optical character verification				
Exper.	Dataset	OK	NOT-OK	Accuracy (%)
1	1	645	645	90.1
2	1	645	444	89.3
	2	2847	2847	96.8
	3	577	577	85.8
2.1	1	714	375	94.8
	2	2954	2954	96.2
	3	199	199	88.1

**Table 4** Experiment results for date character recognition

CNN date character recognition			
Exper.	Dataset	Images per class	Accuracy (%)
3	2	381, 381, 381	92.7
	3	55, 67, 63, 61	90

$$\log \text{softmax}(z_j) = \log \left[ \frac{\exp(z_j)}{\sum_k \exp(z_k)} \right], \quad (31)$$

$z$  is a vector of NN output logits, and  $M$  denotes the batch size of choice for stochastic optimisation of  $\mathcal{L}_{\text{NLL}}$  via backpropagation. In all cases, we use adaptive moment estimate (Adam) as an optimiser [56]. In this framework, three sets of experiments were conducted and the obtained results are reported in Tables 3 and 4.

The goal of the first experiment was to establish a baseline for images that would be classified as acceptable according to human standards. The appearance of unreadable images was especially prominent in the first of the three datasets. Conversely, the average image quality of the second and third datasets was higher; therefore, they were not considered in this experiment. Moreover, the first dataset contained images from seven different locations, and as such, there were at least seven different types of food packaging present. To devise a balanced experiment, images from all locations were combined and categorised into two classes: ‘Complete Dates’ and ‘Unreadable’. As reported in Table 3, **90.1%** classification accuracy was achieved over all seven locations.

The second experiment aimed at distinguishing between acceptable and not-acceptable, missing dates. This meant that the absence of either day or month digits in a *use-by* date is not acceptable. The second dataset was the largest, containing approximately 50% of examples with partial or missing dates. Images missing the day/month or both were assigned to one class and ‘complete dates’ to the other. As reported in Table 3, an accuracy of **96.8%** was achieved.

Similarly, a performance of **94.8%** was achieved when applying the same procedure to the first dataset. As for the third dataset, it includes images of higher quality, but there is a very small number of missing value examples available. To address this, we performed data augmentation in order to produce a larger set of ‘Partial Dates’. The accuracy achieved on this synthetic set was **85.8%**. Lastly, a small variation of this experiment (2.1 in Table 3) was conducted in order to assess how well the network can identify the presence of any type of date, be it complete or partial, versus the absence of a date altogether. This experiment offered insight into how well the network can produce inferred localisation of dates, as it must learn to filter out the abundant non-date-related text/numbers in the images. Table 3 shows that good accuracies were achieved across all three datasets, with the best case of **96.2%** date presence detection on the second dataset.

In a brief third experiment, a global approach to OCV was tested by targeting the classification of specific digits and letters. Successful text recognition systems typically begin with the detection of text presence within a given image, followed by a segmentation or localisation of the desired region-of-interest (ROI) in order to perform classification of segmented digits thereafter. Here, we assess how well the NN can perform without specifying any additional labels or local information. Given that almost all images in the third dataset contained ‘Complete Dates’, we conducted a brief digit classification experiment (see Table 4 for results). Despite the small number of training examples (1138) and limited possible class combinations, four digit classes were identified, namely 5, 8, 16 and 20. With these labelled examples, an accuracy of **90%** was achieved. Similarly for the second dataset—due to limited data—a brief global OCV classification experiment between the months of October and November in *use-by* dates was conducted. An accuracy of **92.7%** was achieved despite the small number of training examples. In reflection of these results, it is important to remember the great variety of text and numbers included in each image. Without providing any local knowledge and given limited training examples, the networks were still able to automatically infer the importance of specific digits and their respective locations in a global manner, whilst ignoring the same or other digits located in close proximity.

## 5.2.2 Latent variable adaptive clustering

A major challenge spanning the three datasets was the high variability in the captured images characteristics. This variability made the reuse of a DNN trained on one dataset, for classifying the data of another, very difficult leading to poor performances. Fundamentally, this is because each dataset comes from a different distribution, as the images

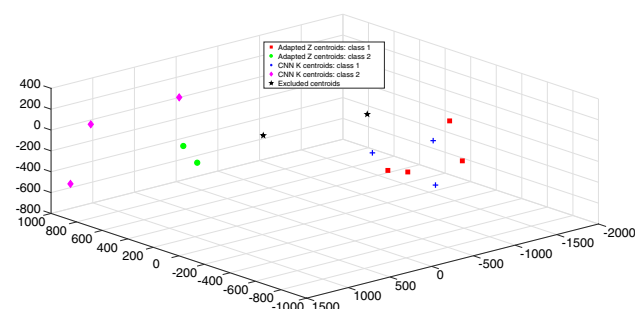


**Table 5** Experiment results of our adaptation procedure

Latent variable adaptive clustering		
Test dataset	Classification accuracy (%)	
	CNN $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$	Our method ( <b>A</b> )
$\mathcal{T}_1$	63.8	76.4
$\mathcal{T}_2$	95.9	97.1

were taken by different people, with different cameras and at differing supplier locations. With limited data available to us, the use of transfer learning among different environments and datasets was ineffective. To overcome these challenges, we demonstrate the possibility of designing a new facet of the same CNN architecture, for learning each considered problem associated with different datasets. The approach focuses on: (i) detecting bad image capturing conditions; (ii) detecting missing dates (i.e. either day and/or month of *use-by* date); (iii) showing the ability to recognise day and/or month of an existing *use-by* date. The CNN architectures proved to be quite accurate in identifying the missing/complete dates classification problem. Subsequently, we explored whether the respective trained networks were suitable for carrying out the proposed network adaptation approach (see Table 5 for results).

To this end, consider  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$  as a trained CNN with a test performance of 95.9% on a binary classification problem of *use-by* date verification on a real dataset. Let  $\mathcal{T}_1$  be the test set of a dataset from a different distribution targeting the same classification task. We forward-propagate  $\mathcal{T}_1$  through  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$  and achieve a lower accuracy of 63.8% as expected. We employed our adaptation procedure to classify  $\mathcal{T}_1$  without any parameter retraining, decreasing the relative error by 34.81% with an improved accuracy of **76.4%**. Interestingly, the original performance achieved by  $\mathcal{F}(\mathcal{D}_2; \mathbf{W}_2)$  on  $\mathcal{T}_2$  also increased from 95.9% to **97.1%** when classifying  $\mathcal{T}_2$  with **A** instead of the CNN, it was originally trained on. Figure 9 depicts a 3D



**Fig. 9** t-SNE visualisation of the derived centroids **A** with best  $k = 7$ , achieving the results reported in Table 5. The ‘Excluded centroids’ (2 black stars) were removed as per the policy outlined in step 6 of our proposed adaptation procedure (colour figure online)

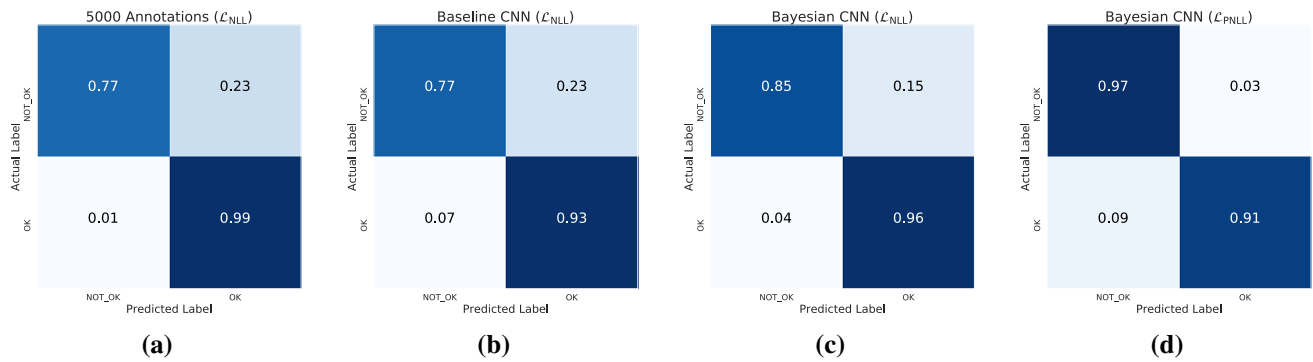
**Table 6** Deep Bayesian self-training performance on real datasets. Cohen’s kappa score  $\kappa$  is also reported

Class	Precision	Recall	F1	#img
Bayesian CNN ( $\mathcal{L}_{\text{PNLL}}$ ), $\kappa = 0.8891$				
NOT-OK	0.9532	0.9694	0.9612	294
OK	0.9427	0.9136	0.9279	162
Avg./total	0.9494	0.9496	0.9494	456
Bayesian CNN ( $\mathcal{L}_{\text{NLL}}$ ), $\kappa = 0.8383$				
NOT-OK	0.9679	0.8538	0.9073	212
OK	0.889	0.9764	0.9306	254
Avg./total	0.9248	0.9206	0.9200	466
Baseline CNN ( $\mathcal{L}_{\text{NLL}}$ ), $\kappa = 0.6964$				
NOT-OK	0.9158	0.7682	0.8355	453
OK	0.7989	0.9287	0.8589	449
Avg./total	0.8576	0.8481	0.8472	902

visualisation of all 2048-dimensional cluster centroids, for  $k = 7$  for both datasets (14 in total). Squares (Red) and (Blue) crosses denote the centroids corresponding to the complete date class in the first and second datasets, respectively. (Green) circles and (Pink) diamonds are the centroids in the missing date category, and the (Black) stars indicate the centroids not used in the final classification as per the centroid exclusion policy explained previously in Sect. 4.

### 5.2.3 Deep Bayesian self-Training on real data

In order to validate our approach, we conducted a series of experiments on a pool of held-out annotated data comprised of 11,948 real food package images. The results can be seen in Table 6 and Fig. 10. We begin by introducing concrete dropout layers after every convolutional layer in the last DenseBlock of a DenseNet-201, pre-trained on ImageNet. We then fine-tuned the last DenseBlock on a small portion of 500 images, with binary annotated labels representing whether the *use-by* date was readable (OK) or not (NOT-OK). As observable in Fig. 10a, we first applied these ideas to the full set of unlabelled 11,948 images and simply selected the 500 most certain predicted labels to be added to the initial training set of 500 images. This process was repeated 10 times in order to collect a total of 5000 images with predicted labels, which we then compared with our annotated labels as shown in Table 6. In the remaining set of experiments, instead of selecting a pre-determined number of images, we filtered out uncertain predictions based on a threshold  $\tau$  as in Algorithm 1. Figure 10c, d depicts the confusion matrices for the automatically annotated images w.r.t. true labels and highlights the benefits of applying a confidence penalty on the log-



**Fig. 10** Normalised confusion matrices of the results obtained from our self-annotation procedure **a** The 5000 predicted labels obtained with the lowest prediction uncertainty. **b** Deterministic baseline CNN predicted labels, wherein the thresholds were set based on the

likelihood loss ( $\mathcal{L}_{PNLL}$ ), as opposed to using a standard log-likelihood ( $\mathcal{L}_{NLL}$ ) which often outputs overconfident distributions. The uncertainties were calculated based on 50 Monte Carlo dropout samples at test time, following the description in Sect. 3.5.

In order to compare our approach to standard self-training, we took the same network and datasets splits and trained it without the Bayesian components. The threshold was set based on the confidence of the CNN output to only consider very confident predictions with over 0.999 predicted probability. As can be seen in Table 6, even with a high threshold, the deterministic CNN tends to be overconfident in its wrong predictions. This causes an increase in the propagated error as more images with wrong predicted labels are added to the training set and the model starts to underperform. To ensure a fair comparison between the self-training methods, the *stop* conditions were set to be identical s.t. the procedure was interrupted after three consecutive iterations without selecting more images to be added to the training set.

## 6 Conclusion and future work

In this paper, we propose a deep Bayesian self-training methodology that leverages modern approximate variational inference in DNNs to estimate predictive uncertainty during a self-training setting. Both aleatoric and epistemic uncertainties of predicted pseudo-labels for unseen data are estimated, and the samples with the lowest predictive uncertainty (highest confidence) are added to the training set in an automated manner. We offer ways to mitigate the known problem of propagating errors in self-training by including: (i) an entropy penalty on the log-likelihood loss to punish overconfident output distributions and facilitate thresholding, and (ii) an adaptive sample-wise weight on the influence of predicted pseudo-labelled samples over

network's sigmoid output. **c** Predicted labels from our Bayesian self-training approach, trained with a standard binary negative log-likelihood loss. **d** Similar to **c** but using a Bayesian CNN trained with the entropy penalised binary negative log-likelihood loss

gradient updates to be inversely proportional to their predictive uncertainty. Lastly, we propose a new simple methodology for visualising and analysing variability between two dataset distributions in DNNs and attempt to adapt information from one problem to the other by clustering learnt latent variable representations in the context of our application domain. An experimental study on both public and private (real) datasets is presented demonstrating the increased performance of our algorithm over standard self-training baselines, and also highlighting the importance of predictive uncertainty estimates in safety-critical domains.

Our future work will extend the experimental study to large dataset, consisting of about half a million real food packaging images, and we intend to apply the presented DNN-based methodologies for adaptation and self-annotation of these data.

**Acknowledgements** The authors would like to thank Mr. George Marandianos, Mrs. Mamatha Thota and Mr. Samuel Bond-Taylor for manually annotating datasets used in this study and of course the reviewers for their constructive feedback that helped to improve the manuscript. We would also like to thank Professor Luc Bidaut for enabling this collaboration.

**Funding** The research presented in this paper was funded by Engineering and Physical Sciences Research Council (Reference Number EP/R005524/1) and Innovate UK (Reference Number 102908), in collaboration with the Olympus Automation Limited Company, for the project Automated Robotic Food Manufacturing System.

## Compliance with ethical standards

**Conflict of Interest** The authors declare that they have no conflict of interest.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a

link to the Creative Commons license, and indicate if changes were made.

## References

- Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*, pp 1097–1105
- Gal Y, Islam R, Ghahramani Z (2017) Deep Bayesian active learning with image data. *arXiv preprint arXiv:1703.02910*
- Zhu X (2006) Semi-supervised learning literature survey. *Comput Sci Univ Wis-Madison* 2(3):4
- Yarowsky D (1995) Unsupervised word sense disambiguation rivaling supervised methods. In: *Proceedings of the 33rd annual meeting on Association for Computational Linguistics*. Association for Computational Linguistics, pp 189–196
- Pratt LY (1993) Discriminability-based transfer between neural networks. In: *Advances in neural information processing systems*, pp 204–211
- Blum A, Mitchell T (1998) Combining labeled and unlabeled data with co-training. In: *Proceedings of the eleventh annual conference on Computational learning theory*. ACM, pp 92–100
- Cohn DA, Ghahramani Z, Jordan MI (1996) Active learning with statistical models. *J Artif Intell Res* 4:129–145
- Zhou Z-H, Li M (2005) Tri-training: exploiting unlabeled data using three classifiers. *IEEE Trans Knowl Data Eng* 11:1529–1541
- Simon T (2001) *Active learning: theory and applications*, vol 1. Stanford University, Stanford
- Han B, Yao Q, Yu X, Niu G, Xu M, Hu W, Tsang I, Sugiyama M (2018) Co-teaching: robust training of deep neural networks with extremely noisy labels. In: *Advances in neural information processing systems*, pp 8527–8537
- Yosinski J, Clune J, Bengio Y, Lipson H (2014) How transferable are features in deep neural networks? In: *Advances in neural information processing systems*, pp 3320–3328
- Zamir AR, Sax A, Shen W, Guibas L, Malik J, Savarese S (2018) Taskonomy: Disentangling task transfer learning. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 3712–3722
- Sousa Ribeiro FD, Calivá F, Swainson M, Gudmundsson K, Leontidis G, Kollias S (2018) An adaptable deep learning system for optical character verification in retail food packaging. In: *IEEE international conference on evolving and adaptive intelligent systems*
- Pearson S, May D, Leontidis G, Swainson M, Brewer S, Bidaut L, Frey JG, Parr G, Maull R (2019) Zisman A (2019) Are distributed ledger technologies the panacea for food traceability? *Glob Food Secur* 20:145–149
- Huang G, Liu Z, Van Der Maaten Laurens, Weinberger KQ (2017) Densely connected convolutional networks. In: *CVPR*, vol 1, p 3
- Sousa Ribeiro FD, Gong L, Calivá F, Swainson M, Gudmundsson K, Yu M, Leontidis G, Ye X, Kollias S (2018) An end-to-end deep neural architecture for optical character verification and recognition in retail food packaging. In: *2018 25th IEEE international conference on image processing (ICIP)*. IEEE, pp 2376–2380
- Sun C, Shrivastava A, Singh S, Gupta A (2017) Revisiting unreasonable effectiveness of data in deep learning era. In: *2017 IEEE international conference on computer vision (ICCV)*. IEEE, pp 843–852
- Chudzik P, Majumdar S, Calivá F, Al-Diri B, Hunter A (2018) Microaneurysm detection using fully convolutional neural networks. *Comput Methods Programs Biomed* 158:185–192
- Kollias D, Yu M, Tagaris A, Leontidis G, Stafylopatis A, Kollias S (2017) Adaptation and contextualization of deep neural network models. In: *2017 IEEE symposium series on computational intelligence (SSCI)*, pp 1–8
- Havaei M, Davy A, Warde-Farley D, Biard A, Courville A, Bengio Y, Pal C, Jodoin P-M, Larochelle H (2017) Brain tumor segmentation with deep neural networks. *Med Image Anal* 35:18–31
- Caliva F, Sousa Ribeiro FD, Mylonakis A, Demaziere C, Vinai P, Leontidis G, Kollias S (2018) A deep learning approach to anomaly detection in nuclear reactors. In: *2018 International joint conference on neural networks (IJCNN)*, pp 1–8
- Sousa Ribeiro FD, Caliva F, Chionis D, Dokhane A, Mylonakis A, Demaziere C, Leontidis G, Kollias S (2018) Towards a deep unified framework for nuclear reactor perturbation analysis. In: *2018 IEEE symposium series on computational intelligence (SSCI)*, pp 1–8
- Kaiser L, Gomez AN, Shazeer N, Vaswani A, Parmar N, Jones L, Uszkoreit J (2017) One model to learn them all. *arXiv preprint arXiv:1706.05137*
- Doersch C, Zisserman A (2017) Multi-task self-supervised visual learning. In: *The IEEE international conference on computer vision (ICCV)*
- Zoph B, Vijay V, Shlens J, Le QV (2017) Learning transferable architectures for scalable image recognition. 2(6). *arXiv preprint arXiv:1707.07012*
- Kingma DP, Mohamed S, Rezende DJ, Welling M (2014) Semi-supervised learning with deep generative models. In: *Advances in neural information processing systems*, pp 3581–3589
- Tzeng E, Hoffman J, Saenko K, Darrell T (2017) Adversarial discriminative domain adaptation. In: *Computer vision and pattern recognition (CVPR)*, vol 1, no 4. pp 7167–7176
- Bousmalis K, Silberman N, Dohan D, Erhan D, Krishnan D (2017) Unsupervised pixel-level domain adaptation with generative adversarial networks. In: *The IEEE conference on computer vision and pattern recognition (CVPR)*, vol 1, no 7
- Tarvainen A, Valpola H (2017) Mean teachers are better role models: weight-averaged consistency targets improve semi-supervised deep learning results. In: *Advances in neural information processing systems*, pp 1195–1204
- Qiao S, Shen W, Zhang Z, Wang B, Yuille A (2018) Deep co-training for semi-supervised image recognition. In: *Proceedings of the European conference on computer vision (ECCV)*, pp 135–152
- Lee D-H (2013) Pseudo-label: the simple and efficient semi-supervised learning method for deep neural networks. In: *Workshop on challenges in representation learning, ICML*, vol 3, pp 2
- Grandvalet Y, Bengio Y (2005) Semi-supervised learning by entropy minimization. In: *Advances in neural information processing systems*, pp 529–536
- Zou Y, Yu Z, Vijaya Kumar BVK, Wang J (2018) Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In: *Proceedings of the European conference on computer vision (ECCV)*, pp 289–305
- Kendall A, Gal Y (2017) What uncertainties do we need in Bayesian deep learning for computer vision? In: *Advances in neural information processing systems*, pp 5574–5584
- Denker JS, Lecun Y (1991) Transforming neural-net output levels to probability distributions. In: *Advances in neural information processing systems*, pp 853–859
- Neal RM (2012) *Bayesian learning for neural networks*, vol 118. Springer, New York

37. MacKay DJC (1992) A practical bayesian framework for back-propagation networks. *Neural Comput* 4(3):448–472
38. Yarín G (2016) Uncertainty in deep learning. University of Cambridge, Cambridge
39. Gal Y, Ghahramani Z (2015) Bayesian convolutional neural networks with Bernoulli approximate variational inference. *arXiv preprint [arXiv:1506.02158](https://arxiv.org/abs/1506.02158)*
40. Kendall A, Gal Y, Cipolla R (2017) Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. *arXiv preprint [arXiv:1705.07115](https://arxiv.org/abs/1705.07115)*
41. Gal Y, Ghahramani Z (2016) Dropout as a Bayesian approximation: representing model uncertainty in deep learning. In: *international conference on machine learning*, pp 1050–1059
42. Hinton GE, Van Camp D (1993) Keeping the neural networks simple by minimizing the description length of the weights. In: *Proceedings of the sixth annual conference on Computational learning theory*, ACM, pp 5–13
43. Graves A (2011) Practical variational inference for neural networks. In: *Advances in neural information processing systems*, pp 2348–2356
44. Welling M, Teh YW (2011) Bayesian learning via stochastic gradient Langevin dynamics. In: *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp 681–688
45. Kingma DP, Welling M (2013) Auto-encoding variational Bayes. *arXiv preprint [arXiv:1312.6114](https://arxiv.org/abs/1312.6114)*
46. Maddison CJ, Mnih A, Teh YW (2016) The concrete distribution: a continuous relaxation of discrete random variables. *arXiv preprint [arXiv:1611.00712](https://arxiv.org/abs/1611.00712)*
47. Gal Y, Hron J, Kendall A (2017) Concrete dropout. In: *Advances in neural information processing systems*, pp 3581–3590
48. Pereyra G, Tucker G, Chorowski J, Kaiser Ł, Hinton G (2017) Regularizing neural networks by penalizing confident output distributions. *arXiv preprint [arXiv:1701.06548](https://arxiv.org/abs/1701.06548)*
49. Kwon Y, Won J-H, Kim BJ, Paik MC (2018) Uncertainty quantification using Bayesian neural networks in classification: application to ischemic stroke lesion segmentation. In: *International conference on medical imaging with deep learning*
50. Arthur D, Vassilvitskii S (2007) k-means++: the advantages of careful seeding. In: *Proceedings of the eighteenth annual ACM-SIAM symposium on discrete algorithms*. Society for Industrial and Applied Mathematics, pp 1027–1035
51. Jégou S, Drozdal M, Vazquez D, Romero A, Bengio Y (2017) The one hundred layers tiramisù: fully convolutional densenets for semantic segmentation. In: *2017 IEEE conference on computer vision and pattern recognition workshops (CVPRW)*. IEEE, pp 1175–1183
52. Ioffe S, Szegedy C (2015) Batch normalization: accelerating deep network training by reducing internal covariate shift. *arXiv preprint [arXiv:1502.03167](https://arxiv.org/abs/1502.03167)*
53. Lakshminarayanan B, Pritzel A, Blundell C (2017) Simple and scalable predictive uncertainty estimation using deep ensembles. In: *Advances in neural information processing systems*, pp 6402–6413
54. Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z (2016) Rethinking the inception architecture for computer vision. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 2818–2826
55. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. *J Mach Learn Res* 15(1):1929–1958
56. Kingma DP, Ba J (2014) Adam: a method for stochastic optimization. *arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980)*

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.